

$$\tilde{Y} = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

Goal: Find estimator  $\hat{\beta}^{(\pi)}$  for some  $\pi$  & compare its risk to  $\hat{\beta}^{(OLS)}$

Bayesian:

- ①  $\pi = N(0, \tau^2 I_p)$   
Compute  $\pi(\beta|Y)$   
 $\hat{\beta}^{(\pi)}$

② Frequentist

③ Compare risk to  $\hat{\beta}^{(OLS)}$

## Linear Regression and Bayesian Estimation

Linear Regression:  $Y_i = X_i^T \beta + \varepsilon_i$ ,  $\varepsilon_i \sim N(0, \sigma^2)$  iid.

$$Y = X\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n), \quad OLS: \hat{\beta}^{(OLS)} = \underset{n \times p \quad n}{\operatorname{argmin}} \|Y - X\beta\|_2^2$$

$$\text{if } X^T X \text{ invertible: } \hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T Y$$

$$\text{quadratic risk: } E[\|\hat{\beta} - \beta\|_2^2]$$

Frequentist: one true (unknown)  $\beta$ , estimate it!  
perform well for any  $\beta$

Bayesian:  $\beta \sim \pi(\cdot)$  prior distribution,  $Y|\beta \sim p(Y|\beta)$

$$\pi(\beta|Y) = \frac{\pi(\beta) \cdot p(Y|\beta)}{\int p(Y|\beta) d\pi(\beta)}$$

$$\hat{\beta}^{(\pi)} = \int \beta d\pi(\beta|Y) \quad \text{posterior mean}$$

$$\vec{y} = \vec{X}\vec{\beta} + \vec{\varepsilon}, \quad \varepsilon \sim N(0, \sigma^2 \mathbf{I}_n)$$

Goal: Find estimator  $\hat{\beta}^{(\pi)}$  for some  $\pi$  & compare its risk to  $\hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T Y$

Bayesian:

$$\textcircled{1} \pi = N(0, \tau^2 \mathbf{I}_p), \lambda = \frac{\sigma^2}{\tau^2}$$

$$\pi(\beta | Y) = N(\mu, \Sigma)$$

$$\Sigma = \sigma^2 (X^T X + \lambda \mathbf{I}_p)^{-1}$$

$$\mu = (X^T X + \lambda \mathbf{I}_p)^{-1} X^T Y = \hat{\beta}^{(\pi)}$$

② Frequentist

③ Compare risk to  $\hat{\beta}^{(OLS)}$

① Posterior distribution

Gaussian density:  $x \sim N(\mu, \Sigma)$

$$\pi(\beta | Y) = \frac{\pi(\beta) \cdot p(Y | \beta)}{p(Y)} \propto \pi(\beta) p(Y | \beta) \left[ \frac{1}{(2\pi)^p \det \Sigma} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)} \right]$$

$$\pi(\beta) p(Y | \beta) = \frac{1}{(2\pi\tau^2)^p} e^{-\frac{1}{2\tau^2} \|\beta\|_2^2} \cdot \frac{1}{(2\pi\sigma^2)^n} e^{-\frac{1}{2\sigma^2} \|Y - X\beta\|_2^2}$$

$$= \frac{1}{(2\pi\tau^2)^p} \cdot \frac{1}{(2\pi\sigma^2)^n} \cdot \exp \left[ -\frac{1}{2\sigma^2} \left( \|Y - X\beta\|_2^2 + \underbrace{\frac{\sigma^2}{\tau^2} \|\beta\|_2^2}_{=\lambda} \right) \right]$$

$$\begin{aligned} \textcircled{*} &= (Y - X\beta)^T (Y - X\beta) + \lambda \beta^T \beta \\ &= Y^T Y - \underbrace{\beta^T X^T Y}_{\sigma^2 \tau^{-1} \frac{1}{\sigma^2} \Sigma X^T Y} - Y^T X \beta + \underbrace{\beta^T X^T X \beta + \lambda \beta^T \beta}_{\beta^T (X^T X + \lambda \mathbf{I}_p) \beta} \end{aligned}$$

$$= \underbrace{\sigma^2 (\beta - \mu)^T \Sigma^{-1} (\beta - \mu)}_{\sigma^2 \mu^T \Sigma \mu} + Y^T Y$$

$$\begin{aligned} \mu &= \frac{1}{\sigma^2} \Sigma X^T Y \\ &= (X^T X + \lambda \mathbf{I}_p)^{-1} X^T Y \end{aligned}$$

$$\vec{y} = \vec{X}\vec{\beta} + \vec{\varepsilon}, \quad \varepsilon \sim N(0, \sigma^2 \mathbf{I}_n)$$

Goal: Find estimator  $\hat{\beta}^{(\pi)}$  for some  $\pi$  & compare its risk to

$$\hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T Y$$

Bayesian:

$$\textcircled{1} \pi = N(0, \tau^2 \mathbf{I}_p), \lambda = \frac{\sigma^2}{\tau^2}$$

$$\pi(\beta | Y) = N(\mu, \Sigma)$$

$$\Sigma = \sigma^2 (X^T X + \lambda \mathbf{I}_p)^{-1}$$

$$\mu = (X^T X + \lambda \mathbf{I}_p)^{-1} X^T Y \\ = \hat{\beta}^{(\pi)}$$

② Frequentist

③ Compare risk to  $\hat{\beta}^{(OLS)}$

②

$$\hat{\beta}^{(R)} = \underset{\beta}{\operatorname{argmin}} \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2$$

$$\text{With } \tau^2 = \frac{\sigma^2}{\lambda}, \hat{\beta}^{(R)} = \hat{\beta}^{(\pi)}$$

"ridge" estimator



$$\tilde{y} = \tilde{X}\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

Goal: Find estimator  $\hat{\beta}^{(\pi)}$  for some  $\pi$  & compare its risk to  $\hat{\beta}^{(OLS)} = (X^T X)^{-1} X^T Y$

Bayesian:

$$\textcircled{1} \pi = N(0, \tau^2 I_p), \lambda = \frac{\sigma^2}{\tau^2}$$

$$\pi(\beta | Y) = N(\mu, \Sigma)$$

$$\Sigma = \sigma^2 (X^T X + \lambda I_p)^{-1}$$

$$\mu = (X^T X + \lambda I_p)^{-1} X^T Y$$

$$= \hat{\beta}^{(\pi)}$$

② Frequentist

③ Compare risk to  $\hat{\beta}^{(OLS)}$

③ Assume  $p=n, X^T X = I_p$ ,  $\hat{\beta}^{(CR)} = (X^T X + \lambda I_p)^{-1} X^T Y$

$$= ((1+\lambda)I_p)^{-1} X^T Y$$

$$= \frac{1}{1+\lambda} X^T Y$$

$$E[\|\hat{\beta}^{(CR)} - \beta\|_2^2]$$

$$= E[\|\frac{1}{1+\lambda} X^T Y - \beta\|_2^2]$$

$$= E[\|\frac{1}{1+\lambda} X^T (X\beta + \varepsilon) - \beta\|_2^2] = E[\|\frac{1}{1+\lambda} \beta - \beta + \frac{1}{1+\lambda} X^T \varepsilon\|_2^2]$$

$$= E[\|\frac{1}{1+\lambda} \beta - \beta\|_2^2] + 2(\frac{1}{1+\lambda} \beta - \beta)^T \underbrace{E[\frac{1}{1+\lambda} X^T \varepsilon]}_{=0}$$

$$+ E[\|\frac{1}{1+\lambda} X^T \varepsilon\|_2^2]$$

$$= \frac{\lambda^2}{(1+\lambda)^2} \|\beta\|_2^2 + \frac{1}{(1+\lambda)^2} E[\varepsilon^T \underbrace{X X^T}_{=I_p} \varepsilon]$$

$$= \frac{\lambda^2}{(1+\lambda)^2} \|\beta\|_2^2 + \frac{p \cdot \sigma^2}{(1+\lambda)^2} = \|\varepsilon\|_2^2 \sim \chi_p^2 \cdot \sigma^2$$

$$\tilde{Y} = \tilde{X}\beta + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I_n)$$

Goal: Find estimator  $\hat{\beta}^{(\pi)}$  for some  $\pi$  & compare its risk to  $\hat{\beta}^{OLS} = (X^T X)^{-1} X^T Y$

Bayesian:  
①  $\pi = N(0, \tau^2 I_p), \lambda = \frac{\sigma^2}{\tau^2}$

$$\begin{aligned} \pi(\beta|Y) &= N(\mu, \Sigma) \\ \Sigma &= \sigma^2 (X^T X + \lambda I_p)^{-1} \\ \mu &= (X^T X + \lambda I_p)^{-1} X^T Y \\ &= \hat{\beta}^{(\pi)} \end{aligned}$$

② Frequentist

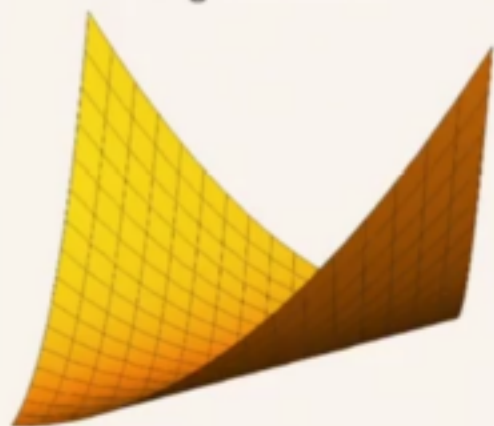
③ Compare risk to  $\hat{\beta}^{OLS}$

$$\textcircled{3} \text{ If } p=n, X^T X = I_p, E[\|\hat{\beta}^{(R)} - \beta\|_2^2] = \underbrace{\frac{\|\beta\|_2^2 \cdot 2^2}{(1+\lambda)^2} + \frac{p\sigma^2}{(1+\lambda)^2}}_{=f(\lambda)}$$

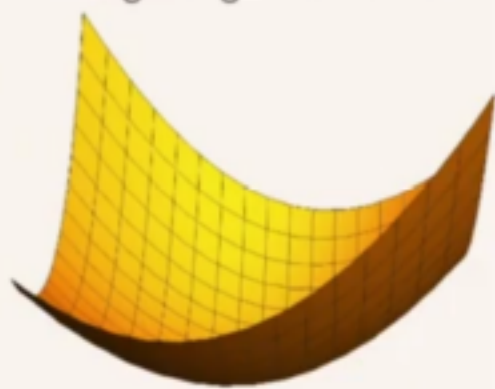
$$\begin{aligned} f'(\lambda) &= \frac{2\lambda \|\beta\|_2^2 \cdot (1+\lambda)^2 - 2^2 \|\beta\|_2^2 \cdot 2(1+\lambda)}{(1+\lambda)^4} \\ &\quad + \frac{-p\sigma^2 \cdot 2(1+\lambda)}{(1+\lambda)^4} \\ &= \frac{2\lambda(1+\lambda)\|\beta\|_2^2 - 2^2 \|\beta\|_2^2 - 2p\sigma^2}{(1+\lambda)^3} \\ &= \frac{2\lambda \|\beta\|_2^2 - 2p\sigma^2}{(1+\lambda)^3} \stackrel{!}{=} 0 \Rightarrow \lambda^* = \frac{p\sigma^2}{\|\beta\|_2^2} \end{aligned}$$

$$\Rightarrow f(\lambda^*) = \frac{\|\beta\|_2^2 \cdot \frac{(p\sigma^2)^2}{\|\beta\|_2^4} + p\sigma^2}{\left(1 + \frac{p\sigma^2}{\|\beta\|_2^2}\right)^2} = p\sigma^2 \cdot \frac{1 + \frac{p\sigma^2}{\|\beta\|_2^2}}{\left(1 + \frac{p\sigma^2}{\|\beta\|_2^2}\right)^2} = p\sigma^2 \cdot \frac{1}{1 + \frac{p\sigma^2}{\|\beta\|_2^2}} < p\sigma^2$$

No regularization



Ridge regularization



$$\hat{\beta}^{(R)} = \underset{\beta}{\operatorname{argmin}} \|X\beta - Y\|_2^2 + \lambda \|\beta\|_2^2$$

Remarks • Obtained estimator that is useful in frequentist setup through Bayes calculation

- $\hat{\beta}^{(R)}$  with "optimal" choice of  $\lambda$  is not a proper estimator:  $\lambda$  depends on  $\beta$ 
  - something known  $\|\beta\|_2^2$
  - cross-validation
  - $X = I_p$ , similar reasoning leads to James-Stein estimator that "beats" OLS
- Especially useful when  $p > n$ ,  $X$  rank deficient